

# UC Berkeley

## UC Berkeley Previously Published Works

**Title**

Data access for the 1,000 Plants (1KP) project.

**Permalink**

<https://escholarship.org/uc/item/1553083t>

**Journal**

GigaScience, 3(1)

**ISSN**

2047-217X

**Authors**

Matasci, Naim  
Hung, Ling-Hong  
Yan, Zhixiang  
[et al.](#)

**Publication Date**

2014

**DOI**

10.1186/2047-217x-3-17

Peer reviewed

REVIEW

Open Access

# Data access for the 1,000 Plants (1KP) project

Naim Matasci<sup>1,2</sup>, Ling-Hong Hung<sup>3</sup>, Zhixiang Yan<sup>4</sup>, Eric J Carpenter<sup>5</sup>, Norman J Wickett<sup>6,7</sup>, Siavash Mirarab<sup>8</sup>, Nam Nguyen<sup>8</sup>, Tandy Warnow<sup>8</sup>, Saravanaraj Ayyampalayam<sup>9</sup>, Michael Barker<sup>2</sup>, J Gordon Burleigh<sup>10</sup>, Matthew A Gitzendanner<sup>10</sup>, Eric Wafula<sup>11</sup>, Joshua P Der<sup>11</sup>, Claude W dePamphilis<sup>11</sup>, Béatrice Roure<sup>12</sup>, Hervé Philippe<sup>12,13</sup>, Brad R Ruhfel<sup>10,14</sup>, Nicholas W Miles<sup>15</sup>, Sean W Graham<sup>16</sup>, Sarah Mathews<sup>17</sup>, Barbara Surek<sup>18</sup>, Michael Melkonian<sup>18</sup>, Douglas E Soltis<sup>10,15,19</sup>, Pamela S Soltis<sup>10,15,19</sup>, Carl Rothfels<sup>20,21</sup>, Lisa Pokorny<sup>20,22</sup>, Jonathan A Shaw<sup>20</sup>, Lisa DeGironimo<sup>23</sup>, Dennis W Stevenson<sup>23</sup>, Juan Carlos Villarreal<sup>24</sup>, Tao Chen<sup>25</sup>, Toni M Kutchan<sup>26</sup>, Megan Rolf<sup>26</sup>, Regina S Baucom<sup>27</sup>, Michael K Deyholos<sup>5</sup>, Ram Samudrala<sup>3</sup>, Zhijian Tian<sup>4</sup>, Xiaolei Wu<sup>4</sup>, Xiao Sun<sup>4</sup>, Yong Zhang<sup>4</sup>, Jun Wang<sup>4</sup>, Jim Leebens-Mack<sup>9\*</sup> and Gane Ka-Shu Wong<sup>4,5,28\*</sup>

## Abstract

The 1,000 plants (1KP) project is an international multi-disciplinary consortium that has generated transcriptome data from over 1,000 plant species, with exemplars for all of the major lineages across the *Viridiplantae* (green plants) clade. Here, we describe how to access the data used in a phylogenomics analysis of the first 85 species, and how to visualize our gene and species trees. Users can develop computational pipelines to analyse these data, in conjunction with data of their own that they can upload. Computationally estimated protein-protein interactions and biochemical pathways can be visualized at another site. Finally, we comment on our future plans and how they fit within this scalable system for the dissemination, visualization, and analysis of large multi-species data sets.

**Keywords:** *Viridiplantae*, Biodiversity, Transcriptomes, Phylogenomics, Interactions, Pathways

## Introduction

The 1,000 plants (1KP) project is an international multi-disciplinary consortium that has now generated transcriptome data from over 1,000 plant species. One of the goals of our species selection process was to provide exemplars for all of the major lineages across the *Viridiplantae* (green plants), representing approximately one billion years of evolution, including flowering plants, conifers, ferns, mosses and streptophyte green algae. Whereas genomics has long strived for completeness within species (e.g., every gene in the species), we were focused on completeness across an evolutionary clade – obviously not every species, but one representative species for everything at some phylogenetic level (e.g., one species per family, and perhaps more than one species when the family is especially large). Because many of our species had never been subjected to large-scale sequencing, 2 gigabases (Gb) of data per sample was sufficient to increase the number of plant genes by

approximately 100-fold in comparison to the totality of the public databases.

The 1KP project began as a public-private partnership, with 75% of the funding provided by the Government of Alberta and 25% by Musea Ventures. Significant in-kind contributions were provided by BGI-Shenzhen in the form of reduced sequencing costs and by the NSF-funded iPlant collaborative [1] in the form of computational informatics support. Many plant scientists from around the world were involved in the collection of live tissue samples and in the extraction of RNA. Additional computing resources were provided by Compute Canada and by the China National GeneBank. Despite the constraints of this funding model, we released our data (on a collaborative basis) to scientists who approached us with goals that did not compete with ours. For the general community, access was provided through a BLAST portal [2].

We believed that there would be intrinsic value in data of this nature that is beyond our imagination. But for the initial publication, we agreed on two objectives. Firstly, by adopting a phylogenomics approach we hoped to resolve many of the lingering uncertainties in species

\* Correspondence: jleebensmack@plantbio.uga.edu; gane@ualberta.ca

<sup>9</sup>Department of Plant Biology, University of Georgia, Athens, GA, 30602, USA

<sup>4</sup>BGI-Shenzhen, Bei Shan Industrial Zone, Shenzhen, China

Full list of author information is available at the end of the article

relationships, especially in the early lineages of streptophyte green algae and land plants, where previous analyses were based on comparatively sparse taxonomic densities. And secondly, despite the limitations of these data, we hoped to identify some of the gene changes associated with the major innovations in *Viridiplantae* evolution, such as multicellularity, transitions from marine to freshwater or terrestrial environments, maternal retention of zygotes and embryos, complex life history involving haploid and diploid phases, vascular systems, seeds and flowers.

Our RNA extraction protocols [3] and our RNA-Seq transcriptome assembly algorithms [4] have already been published. Here, we are publishing the second of two linked papers. The first is a review of the state-of-knowledge for *Viridiplantae* species relationships and our initial foray into the phylogenomics on a subset of 1KP [5]. The other is a description of the websites that we created in order to provide access to the data (from raw reads to computed results), visualize the results, and perform custom analyses in conjunction with external data that the users can upload. An initial gene annotation is also provided, which focuses on the functional relationships between proteins and their associated metabolites.

## Review

### Access to raw and processed data

Our initial phylogenomics effort used sequences from multiple sources. They include transcriptomes from 1KP representing 85 species, transcriptomes from other sources representing 7 species, and genomes representing an additional 11 species. A summary of these data sources is given in Table 1. We submitted all of the unassembled reads from the 1KP transcriptomes to the Short Reads Archive (SRA) under project accession PRJEB4921 “1000 Plant (1KP) Transcriptome: The Pilot Study.” Note that, with the exception of *Eschscholzia californica*, we sequenced only one sample per species.

To make it easier for others to reproduce our phylogenomics analyses, we are releasing our intermediate computations, not just the final results. Everything is hosted at the iPlant Data Store, a high performance, large capacity, distributed storage system. The contents include transcriptome assemblies, putative coding sequences, orthogroups (i.e., from the 11 reference genomes), as well as gene and species trees with related sequence alignments. There are quite a lot of files and their total sizes are not negligible; so before users begin to download these files, we suggest that they consult Table 2 for a description of what to expect.

At the simplest level, anonymous downloads are permitted from a designated area of the iPlant Data Store [6]. However, much greater functionality is available through the iPlant resources that we describe in the following sections.

### Visualization and custom analyses

To take full advantage of the iPlant computational infrastructure, it is necessary to first register at [7]. Accounts are free, and in addition to 1KP data, users will find high performance computing and cloud-based services. Multiple access modalities are supported: anonymous and secure web interfaces, desktop clients and high-speed command lines. However, we feel that for most users the best option is the iPlant discovery environment (DE), a web-based interface that provides users with high-performance computing resources and data storage. Most contemporary web browsers are supported, including Safari v. 6.1, Firefox v. 24, and Chrome v. 34. The caveat is that some of these functionalities (see below) require Java 1.6.

To guide users through its resources, iPlant is constantly producing new tutorials and teaching materials, including live and recorded webinars. The full catalog can be found at [8]. Here, we describe the new resources specifically created for 1KP.

### Discovery environment (DE)

For access to the 1KP files, users should visit [9] and search for a folder called *Community Data/onekp\_pilot* Figure 1.

From the data window it is possible to download individual files or perform bulk downloads of multiple files and directories through a Java plugin. Note that for security reasons, some operating systems will not allow users to run Java applets. In this instance, a window will pop up to tell the user that there is a problem, and the user should follow the instructions that are given to configure an iDrop desktop [10] Figure 2.

It is possible to perform analyses directly in the DE using any of the 1KP files as input; for example, users can re-compute the sequence alignments and gene trees using different algorithms and parameters [11] Figure 3. More generally, users can select from a variety of applications in the Apps catalogue, which is constantly growing, and includes many popular bioinformatics tools for large-scale phylogenetics, genome-wide associations and next generation sequence analyses.

Species and gene trees can be explored with the iPlant tree viewer, *Phylozoom*, a newly developed web-based phylogenetic tree viewer that supports trees with hundreds of thousand leaves and allows for semantic zooming Figure 4. To access the tree viewer, users need only click on a tree file. This will open a preview window with two tabs: one for the tree's newick string (a format for graph-theoretical trees as defined at [12]) and another for the web link that opens a window to the tree display. Notice that pop-ups must be enabled on the user's browser.

To zoom in and expand the collapsed clades, click on the node of interest. To zoom out, click and drag the tree

**Table 1 Data sources for phylogenomics analyses**

Species	Type	Accession	iPlant ID
<i>Arabidopsis thaliana</i>	genome	n/a	n/a
<i>Brachypodium distachyon</i>	genome	n/a	n/a
<i>Carica papaya</i>	genome	n/a	n/a
<i>Medicago truncatula</i>	genome	n/a	n/a
<i>Oryza sativa</i>	genome	n/a	n/a
<i>Physcomitrella patens</i>	genome	n/a	n/a
<i>Populus trichocarpa</i>	genome	n/a	n/a
<i>Selaginella moellendorffii</i>	genome	n/a	n/a
<i>Sorghum bicolor</i>	genome	n/a	n/a
<i>Vitis vinifera</i>	genome	n/a	n/a
<i>Zea mays</i>	genome	n/a	n/a
<i>Aquilegia formosa</i>	meta-assembly	PlantGDB	AQUI
<i>Cycas rumphii</i>	meta-assembly	SRX022306, SRX022215	CYCA
<i>Liriodendron tulipifera</i>	meta-assembly	PRJNA46857	LIRI
<i>Persea americana</i>	meta-assembly	PRJNA46857	PERS
<i>Pinus taeda</i>	meta-assembly	PRJNA79733	PINU
<i>Pteridium aquilinum</i>	meta-assembly	PRJNA48473	PTER
<i>Zamia vazezizii</i>	meta-assembly	PRJNA46857	ZAMI
<i>Acorus americanus</i>	OneKP meta-assembly	ERR364395, PRJNA46857	ACOR
<i>Amborella trichopoda</i>	OneKP meta-assembly	ERR364329, PRJNA46857	AMBO
<i>Catharanthus roseus</i>	OneKP meta-assembly	ERR364390, PRJNA79951, PRJNA236160	CATH
<i>Eschscholzia californica</i>	OneKP meta-assembly	ERR364338, ERR364335, ERR364336, ERR364337, ERR364334, SRX002988, SRX002987, PlantGDB	ESCH
<i>Ginkgo biloba</i>	OneKP meta-assembly	ERR364401, PlantGDB	GINK
<i>Nuphar advena</i>	OneKP meta-assembly	ERR364330, PRJNA46857	NUPH
<i>Ophioglossum petiolatum</i>	OneKP meta-assembly	ERR364410, SRX666586	OPHI
<i>Saruma henryi</i>	OneKP meta-assembly	ERR364383, PRJNA46857	SARU
<i>Welwitschia mirabilis</i>	OneKP meta-assembly	ERR364404, PRJNA46857	WELW
<i>Allamanda cathartica</i>	OneKP	ERR364389	MGVU
<i>Angiopteris evecta</i>	OneKP	ERR364409	NHCM
<i>Anomodon attenuatus</i>	OneKP	ERR364349	QMWB
<i>Bazzania trilobata</i>	OneKP	ERR364415	WZYK
<i>Boehmeria nivea</i>	OneKP	ERR364387	ACFP
<i>Bryum argenteum</i>	OneKP	ERR364348	JMXW
<i>Cedrus libani</i>	OneKP	ERR364342	GGEA
<i>Ceratodon purpureus</i>	OneKP	ERR364350	FFPD
<i>Chaetosphaeridium globosum</i>	OneKP	ERR364369	DRGY
<i>Chara vulgaris</i>	OneKP	ERR364366	CHAR
<i>Chlorokybus atmophyticus</i>	OneKP	ERR364371	AZZW
<i>Colchicum autumnale</i>	OneKP	ERR364397	NHIX
<i>Coleochaete irregularis</i>	OneKP	ERR364367	QPDY
<i>Coleochaete scutata</i>	OneKP	ERR364368	VQBJ
<i>Cosmarium ochthodes</i>	OneKP	ERR364376	STKJ
<i>Cunninghamia lanceolata</i>	OneKP	ERR364340	OUOI

**Table 1 Data sources for phylogenomics analyses (Continued)**

<i>Cyathea (Alsophila) spinulosa</i>	OneKP	ERR364412	GANB
<i>Cycas micholitzii</i>	OneKP	ERR364405	XZUY
<i>Cylindrocystis brebissonii</i>	OneKP	ERR364378	YOXI
<i>Cylindrocystis cushleackae</i>	OneKP	ERR364373	JOJQ
<i>Dendrolycopodium obscurum</i>	OneKP	ERR364346	XNXF
<i>Dioscorea villosa</i>	OneKP	ERR364396	OCWZ
<i>Diospyros malabarica</i>	OneKP	ERR364339	KVFU
<i>Entransia fimbriata</i>	OneKP	ERR364372	BFIK
<i>Ephedra sinica</i>	OneKP	ERR364402	VDAO
<i>Equisetum diffusum</i>	OneKP	ERR364408	CAPN
<i>Gnetum montanum</i>	OneKP	ERR364403	GTHK
<i>Hedwigia ciliata</i>	OneKP	ERR364352	YWNF
<i>Hibiscus cannabinus</i>	OneKP	ERR364388	OLXF
<i>Houttuynia cordata</i>	OneKP	ERR364332	CSSK
<i>Huperzia squarrosa</i>	OneKP	ERR364407	GAON
<i>Inula helenium</i>	OneKP	ERR364393	AFQQ
<i>Ipomoea purpurea</i>	OneKP	ERR364392	VXKB
<i>Juniperus scopulorum</i>	OneKP	ERR364341	XMGP
<i>Kadsura heteroclita</i>	OneKP	ERR364331	NWMY
<i>Klebsormidium subtile</i>	OneKP	ERR364370	FQLP
<i>Kochia scoparia</i>	OneKP	ERR364385	WGET
<i>Larrea tridentata</i>	OneKP	ERR364386	UDUT
<i>Leucodon brachypus</i>	OneKP	ERR364353	ZACW
<i>Marchantia emarginata</i>	OneKP	ERR364417	TFYI
<i>Marchantia polymorpha</i>	OneKP	ERR364416	JPYU
<i>Mesostigma viride</i>	OneKP	ERR364365	KYIO
<i>Mesotaenium endlicherianum</i>	OneKP	ERR364377	WDCW
<i>Metzgeria crassipilis</i>	OneKP	ERR364359	NRWZ
<i>Monomastix opisthostigma</i>	OneKP	ERR364362	BTfM
<i>Mougeotia</i> sp.	OneKP	ERR364374	ZRMT
<i>Nephroselmis pyriformis</i>	OneKP	ERR364363	ISIM
<i>Netrium digitus</i>	OneKP	ERR364379	FFGR
<i>Nothoceros aenigmaticus</i>	OneKP	ERR364356	DXOU
<i>Nothoceros vincentianus</i>	OneKP	ERR364357	TCBC
<i>Penium margaritaceum</i>	OneKP	ERR364382	AEKF
<i>Podophyllum peltatum</i>	OneKP	ERR364384	WFBF
<i>Polytrichum commune</i>	OneKP	ERR364413	SZYG
<i>Prumnopitys andina</i>	OneKP	ERR364343	EGLZ
<i>Pseudolycopodiella caroliniana</i>	OneKP	ERR364345	UPMJ
<i>Psilotum nudum</i>	OneKP	ERR364411	QVMR
<i>Pyramimonas parkeae</i>	OneKP	ERR364361	TNAW
<i>Rhynchostegium serrulatum</i>	OneKP	ERR364355	JADL
<i>Ricciocarpos natans</i>	OneKP	ERR364358	WJLO
<i>Rosmarinus officinalis</i>	OneKP	ERR364391	FDMM
<i>Rosulabryum</i> cf. <i>capillare</i>	OneKP	ERR364351	XWHK

**Table 1 Data sources for phylogenomics analyses (Continued)**

<i>Roya obtusa</i>	OneKP	ERR364380	XRTZ
<i>Sabal bermudana</i>	OneKP	ERR364400	HWUP
<i>Sarcandra glabra</i>	OneKP	ERR364333	OSHQ
<i>Sciadopitys verticillata</i>	OneKP	ERR364344	YFZK
<i>Selaginella stauntoniana</i>	OneKP	ERR364347	ZZOL
<i>Smilax bona-nox</i>	OneKP	ERR364398	MWYQ
<i>Sphaerocarpos texanus</i>	OneKP	ERR364360	HERT
<i>Sphagnum lescurii</i>	OneKP	ERR364414	GOWD
<i>Spirogyra</i> sp.	OneKP	ERR364375	HAOX
<i>Spirotaenia minuta</i>	OneKP	ERR364381	NNHQ
<i>Tanacetum parthenium</i>	OneKP	ERR364394	DUQG
<i>Taxus baccata</i>	OneKP	ERR364406	WWSS
<i>Thuidium delicatulum</i>	OneKP	ERR364354	EEMJ
<i>Uronema</i> sp.	OneKP	ERR364364	ISGT
<i>Yucca filamentosa</i>	OneKP	ERR364399	ICNN

Meta-assembly refers to a transcriptome assembled from more than one sequenced sample. Some of these were a combination of 1KP and other data; some were entirely non-1KP. Accession numbers (SRA or otherwise) are given for all of the transcriptomes that we used.

figure to the left. To zoom out completely, click the space bar. The web address is a unique identifier that can be shared with others to let them to visualize the tree.

For more advanced users wanting to perform more complicated procedures, iPlant capabilities are available from a command line. It is based on the integrated rule-oriented data system (iRODS) [13]. All the user has to do is install a command line utility, *icommands*, which mimics UNIX and enables high-speed parallel data transfers. Instructions are available at [14].

### Interactions and pathways

In addition to the tree-based species and gene relationships at the iPlant site, functional relationships between proteins and their associated metabolites are available from the Computational Biology Group at the University of Washington, developers of CANDO [15]. Sequence similarity-based methods are used to map 1KP proteins to curated repositories of protein-protein interactions (i.e., BioGRID [16]) and biochemical pathways (i.e., Kyoto Encyclopedia of Genes and Genomes [KEGG] [17]). The user can select any metabolic pathway defined by KEGG and, within this context, see all the 1KP proteins from their chosen species, with functional annotations inferred from KEGG. This website is at [18] Figure 5.

Note that, over the course of this project, there have been many improvements in the transcriptome assemblies. The phylogenomics work (now being published) was done with the SOAPdenovo algorithm. A second assembly was subsequently done with the

newer SOAPdenovo-trans algorithm, which we incorporated into the newer interactions and pathways work. However, both sets of assemblies are available through the iPlant data store.

### Conclusions

The rest of the 1KP data will be released, on much the same platform, along with our analyses of all one thousand species. Our scientific objectives are given at [19]. We have always been open about our intentions, because we wanted to avoid conflict among the scientists who were already working with 1KP and offer early pre-publication access to other non-competing scientists. As soon as we see a draft of a paper, we track its progress through the review process at [20]. Some of these papers have already been published, and more than a few required years of follow-up experiments, resulting for example in fundamental discoveries for molecular evolution [21] and (surprisingly) new tools for mammalian neurosciences [22].

Many of these studies were not anticipated when 1KP was conceived. We only knew that, just as there was value in sequencing every gene in a genome, despite not knowing *a priori* what many of the genes might do, there would be value in sequencing across an ancient and ecologically dominant clade, even when many of the species have no obvious economic or scientific value that would justify a genome sequencing effort. Transcriptomes were a less expensive way to explore plant diversity, and demonstrate value beyond the obvious species.

**Table 2 Number and size of data files on websites**

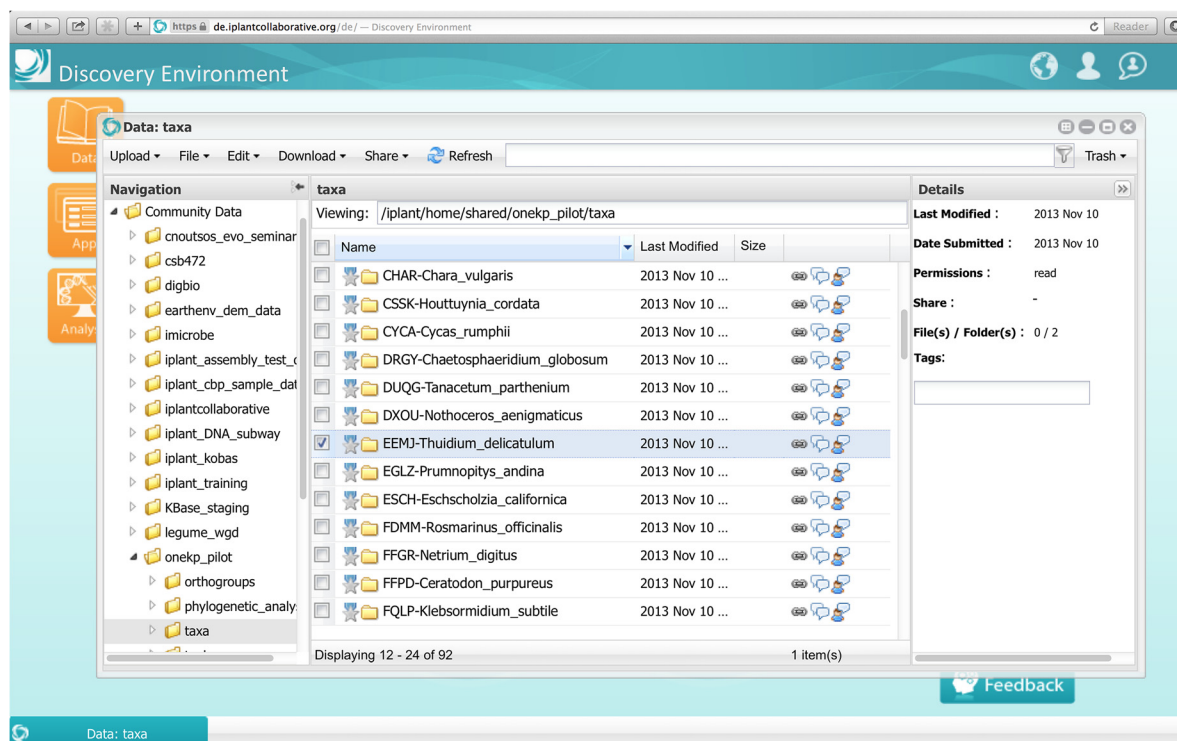
File count	Median size (Mb)	Average size (Mb)	Largest size (Mb)	Total size (Mb)	Similar directories	iPlant directory name
68,253	0.0	0.3	481.1	23,116.6		onekp_pilot
48,053	0.0	0.3	481.1	14,956.7		onekp_pilot/orthogroups
19,220	0.1	0.7	243.8	13,276.5		onekp_pilot/orthogroups/alignments
9,610	0.1	0.3	79.8	3,289.6		onekp_pilot/orthogroups/alignments/FAA
9,610	0.2	1.0	243.8	9,986.9		onekp_pilot/orthogroups/alignments/FNA
28,833	0.0	0.1	481.1	1,680.2		onekp_pilot/orthogroups/gene_trees
9,611	0.0	0.1	481.1	583.3		onekp_pilot/orthogroups/gene_trees/FAA
9,610	0.0	0.0	0.5	102.2		onekp_pilot/orthogroups/gene_trees/FAA/trees
19,222	0.0	0.1	458.0	1,096.8		onekp_pilot/orthogroups/gene_trees/FNA
9,611	0.0	0.1	458.0	556.6		onekp_pilot/orthogroups/gene_trees/FNA/12_codon
9,610	0.0	0.0	0.5	98.5		onekp_pilot/orthogroups/gene_trees/FNA/12_codon/trees
9,611	0.0	0.1	439.1	540.3		onekp_pilot/orthogroups/gene_trees/FNA/all_codon
9,610	0.0	0.0	0.5	101.2		onekp_pilot/orthogroups/gene_trees/FNA/all_codon/dna_tree
19,919	0.0	0.2	175.2	3,468.8		onekp_pilot/phylogenetic_analysis
2,556	0.1	0.1	1.0	292.7		onekp_pilot/phylogenetic_analysis/alignments
852	0.0	0.0	0.3	41.8		onekp_pilot/phylogenetic_analysis/alignments/FAA
852	0.1	0.1	1.0	125.5		onekp_pilot/phylogenetic_analysis/alignments/FNA
852	0.1	0.1	0.9	125.4		onekp_pilot/phylogenetic_analysis/alignments/FNA2AA
17,197	0.0	0.1	0.4	1,827.3		onekp_pilot/phylogenetic_analysis/gene_trees
1,704	0.0	0.1	0.4	238.3		onekp_pilot/phylogenetic_analysis/gene_trees/FAA
2	0.3	0.1	0.4	0.3	852	onekp_pilot/phylogenetic_analysis/gene_trees/FAA/raxmlboot.####
1,704	0.0	0.1	0.4	238.3		onekp_pilot/phylogenetic_analysis/gene_trees/FNA
2	0.3	0.1	0.4	0.3	852	onekp_pilot/phylogenetic_analysis/gene_trees/FNA/raxmlboot.####
3,408	0.0	0.1	0.4	476.7		onekp_pilot/phylogenetic_analysis/gene_trees/FNA2AA
2	0.3	0.1	0.4	0.3	852	onekp_pilot/phylogenetic_analysis/gene_trees/FNA2AA/raxmlboot.####
2	0.3	0.1	0.4	0.3	852	onekp_pilot/phylogenetic_analysis/gene_trees/FNA2AA/raxmlboot.####.c1c2
10,381	0.0	0.1	0.4	874.0		onekp_pilot/phylogenetic_analysis/gene_trees/filtered
2,548	0.0	0.1	0.4	169.3		onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FAA
1	0.0	0.0	0.0	0.0	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FAA/raxmlboot.####.f25
1	0.2	0.1	0.4	0.2	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FAA/raxmlboot.####.filterlen33
852	0.0	0.0	0.0	3.8		onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA
1	0.0	0.0	0.0	0.0	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA/raxmlboot.####.f25
6,980	0.0	0.1	0.4	700.9		onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA
2	0.3	0.1	0.4	0.3	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA/raxmlboot.####.GAMMA.2
2	0.3	0.1	0.4	0.3	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA/raxmlboot.####.c1c2.GAMMA.2
1	0.0	0.0	0.0	0.0	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA/raxmlboot.####.c1c2.f25



**Table 2 Number and size of data files on websites (Continued)**

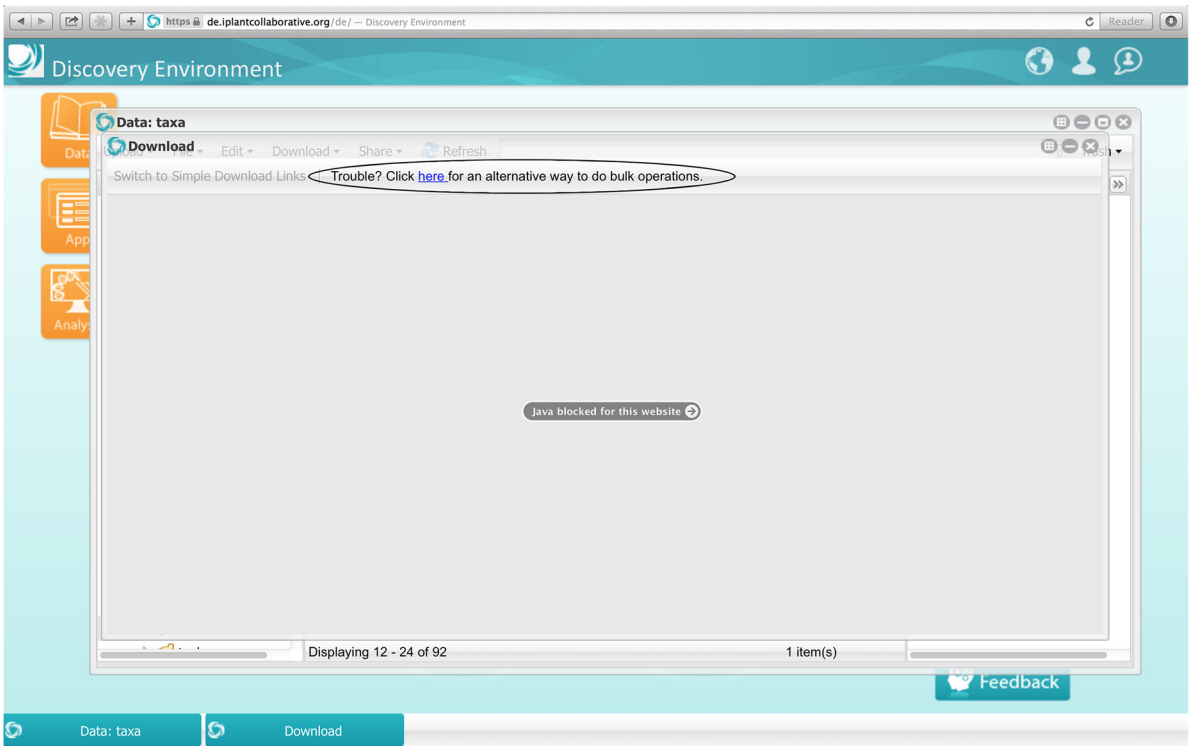
1	0.0	0.0	0.0	0.0	852	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA/raxmlboot.####.f25
2	0.2	0.1	0.4	0.2	844	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA/raxmlboot.####.filterlen33
1	0.3	0.3	0.4	0.3	180	onekp_pilot/phylogenetic_analysis/gene_trees/filtered/FNA2AA/raxmlboot.####.filtered25.GAMMA.2
166	0.0	8.1	175.2	1,348.8		onekp_pilot/phylogenetic_analysis/species_level
50	15.0	27.0	175.2	1,348.1		onekp_pilot/phylogenetic_analysis/species_level/alignments
15	14.7	14.3	58.3	214.2		onekp_pilot/phylogenetic_analysis/species_level/alignments/FAA
35	29.4	32.4	175.2	1,133.9		onekp_pilot/phylogenetic_analysis/species_level/alignments/FNA
116	0.0	0.0	0.0	0.6		onekp_pilot/phylogenetic_analysis/species_level/trees
276	10.0	17.0	157.4	4,691.1		onekp_pilot/taxa
3	9.7	17.0	157.4	51.0	92	onekp_pilot/taxa/####-#####
1	30.8	17.0	157.4	36.0	92	onekp_pilot/taxa/####-#####/assemblies
2	9.7	7.5	45.2	15.0	92	onekp_pilot/taxa/####-#####/translations
5	0.0	0.0	0.0	0.1		onekp_pilot/tools
File count	Median size (Mb)	Average size (Mb)	Largest size (Mb)	Total size (Mb)	Similar directories	Contents at SRA (PRJEB4921)
178	1,915.0	2,045.5	3,371.0	364,100.0		total of all short reads – uncompressed, but downloads are compressed to a quarter of these sizes
2	1,915.0	2,045.5	3,371.0	4,091.0	89	expecting per sample – uncompressed, but downloads are compressed to a quarter of these sizes

In some instances, users will find many directories with similar names, as indicated in this table by hash (#) marks. The total number of directories is given in the preceding column.

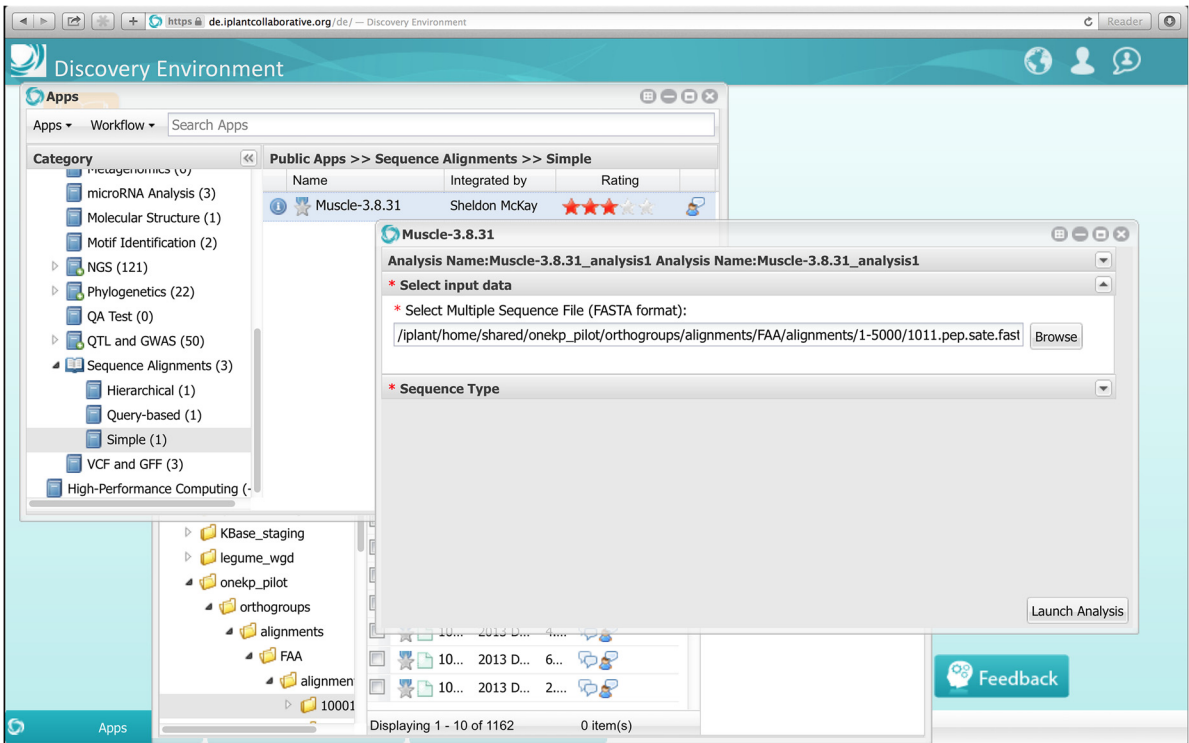


**Figure 1 iPlant DE data window.**

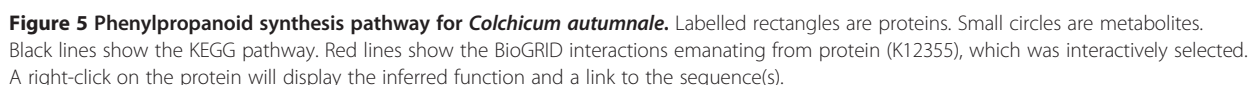




**Figure 2** Bulk download window if Java is disabled. Click on the circled link to access the instructions to install and configure an iDrop desktop.



**Figure 3** Realigning a group of sequences using *Muscle*.



## Abbreviations

1KP: 1,000 Plants project; DE: Discovery Environment; KEGG: Kyoto Encyclopedia of Genes and Genomes; NSF: National Science Foundation; SRA: Short Reads Archive.

## Competing interests

The authors declare that they have no competing interests.

## Authors' contributions

CWD, BRR, NWM, SWG, S Ma, BS, MM, DES, PSS, CR, LP, JAS, LD, DWS, JCV, TC, TMK, MR, RSB, MKD, and JLM collected the plant samples. NM, NJW, S Mi, NN, TW, SA, MB, JGB, MAG, EW, JPD, CWD, BR, HP, BRR, and JLM performed the phylogenomic analyses. NM, LHH, ZY, and EJC setup and maintained web-resources used to communicate data. LHH and RS performed the protein and KEGG pathway analyses. EJC, ZT, XW, XS, YZ, JW, and GKW generated the sequence data. GKW and JLM designed and oversaw the research. All authors read and approved the final manuscript.

## Acknowledgments

The 1000 Plants (1KP) initiative, led by GKW, is funded by the Alberta Ministry of Innovation and Advanced Education, Alberta Innovates Technology Futures (AITF), Innovates Centre of Research Excellence (iCORE), Musea Ventures, BGI-Shenzhen and China National GeneBank (CNCB). We thank the many people responsible for sample collection on 1KP and the staff at BGI-Shenzhen for doing our sequencing. Phylogenomic analyses were supported by the US National Science Foundation through the iPlant collaborative. CANDO was funded by an NIH Director's Pioneer Award 1DP1OD006779-01.

## Author details

<sup>1</sup>iPlant Collaborative, Tucson 85721, AZ, USA. <sup>2</sup>Department of Ecology and Evolutionary Biology, University of Arizona, Tucson 85721, AZ, USA. <sup>3</sup>Department of Microbiology, University of Washington, Seattle 98109, WA, USA. <sup>4</sup>BGI-Shenzhen, Bei Shan Industrial Zone, Shenzhen, China. <sup>5</sup>Department of Biological Sciences, University of Alberta, Edmonton T6G 2E9, AB, Canada. <sup>6</sup>Chicago Botanic Garden, Glencoe 60022, IL, USA. <sup>7</sup>Program in Biological Sciences, Northwestern University, Evanston 60208, IL, USA. <sup>8</sup>Department of Computer Science, University of Texas, Austin, TX, 78712, USA. <sup>9</sup>Department of Plant Biology, University of Georgia, Athens, GA, 30602, USA. <sup>10</sup>Department of Biology, University of Florida, Gainesville, FL 32611, USA. <sup>11</sup>Department of Biology, Penn State University, University Park, Pennsylvania, PA, 16801, USA. <sup>12</sup>Département de Biochimie, Centre Robert-Cedergren, Université de Montréal, Succursale Centre-Ville, Montréal, Québec H3C3J7, Canada. <sup>13</sup>CNRS, USR 2936, Station d'Ecologie Expérimentale du CNRS, Moulis 09200, France. <sup>14</sup>Department of Biological Sciences, Eastern Kentucky University, Richmond, KY, 40475, USA. <sup>15</sup>Florida Museum of Natural History, Gainesville, FL, 32611, USA. <sup>16</sup>Department of Botany, University of British Columbia, Vancouver, BC V6T 1Z4, Canada. <sup>17</sup>Arnold Arboretum of Harvard University, Cambridge, MA, 02138, USA. <sup>18</sup>Botanical Institute, Universität zu Köln, Köln D-50674, Germany. <sup>19</sup>Genetics Institute, University of Florida, Gainesville, FL, 32611, USA. <sup>20</sup>Department of Biology, Duke University, Durham, NC 27708, USA. <sup>21</sup>Department of Zoology, University of British Columbia, Vancouver, BC, V6T 1Z4, Canada. <sup>22</sup>Department of Biodiversity and Conservation, Real Jardín Botánico (RJB-CSIC), 28014 Madrid, Spain. <sup>23</sup>New York Botanical Garden, Bronx, NY, 10458, USA. <sup>24</sup>Systematic Botany and Mycology, University of Munich (LMU), Menzinger Str. 67, 80638 Munich, Germany. <sup>25</sup>Shenzhen Fairy Lake Botanical Garden, The Chinese Academy of Sciences, Shenzhen, Guangdong, 518004, China. <sup>26</sup>Donald Danforth Plant Science Center, St. Louis, MO, 63132, USA. <sup>27</sup>Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, MI, USA. <sup>28</sup>Department of Medicine, University of Alberta, Edmonton, AB, T6G 2E1, Canada.

Received: 22 May 2014 Accepted: 2 October 2014

Published: 27 October 2014

## References

- Goff SA, Vaughn M, McKay S, Lyons E, Stapleton AE, Gessler D, Matasci N, Wang L, Hanlon M, Lenards A, Muir A, Merchant N, Lowry S, Mock S, Helmke M, Kubach A, Narro M, Hopkins N, Micklos D, Hilgert U, Gonzales M, Jordan C, Skidmore E, Dooley R, Cazes J, McLay R, Lu Z, Pasternak S, Koesterke L, Piel WH, et al: **The iPlant Collaborative: Cyberinfrastructure for Plant Biology.** *Front Plant Sci* 2011, **2**:34.

- 1KP BLAST Search Portal** [http://www.bioinfodata.org/app/Blast4OneKP/home]
- Johnson MT, Carpenter EJ, Tian Z, Bruskewich R, Burris JN, Carrigan CT, Chase MW, Clarke ND, Covshoff S, dePamphilis CW, Edger PP, Goh F, Graham S, Greiner S, Hibberd JM, Jordon-Thaden I, Kutchan TM, Leebens-Mack J, Melkonian M, Miles N, Myburg H, Patterson J, Pires JC, Ralph P, Rolf M, Sage RF, Soltis D, Soltis P, Stevenson D, Stewart CN Jr, et al: **Evaluating methods for isolating total RNA and predicting the success of sequencing phylogenetically diverse plant transcriptomes.** *PLoS One* 2012, **7**:e50226.
- Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Li S, Zhou X, Lam TW, Li Y, Xu X, Wong GK, Wang J: **SOAPdenovo-Trans: De novo transcriptome assembly with short RNA-Seq reads.** *Bioinformatics* 2014, **30**:1660–1666.
- Wickett NJ, Mirarab S, Nguyen N, Warnow T, Carpenter E, Matasci N, Ayyampalayam S, Barker M, Burleigh JG, Gitzendanner MA, Ruhfel B, Wafula E, Der JP, Graham SW, Mathews S, Melkonian M, Soltis DE, Soltis PS, Miles NW, Rothfels C, Pokorny L, Shaw AJ, deGironimo L, Stevenson DW, Surek B, Villarreal JC, Roure B, Philippe H, dePamphilis CW, Chen T, et al: **A phylotranscriptomics analysis of the origin and early diversification of land plants.** *Proc Natl Acad Sci U S A* IN PRESS.
- iPlant Data Store for 1KP Pilot.** [http://mirrors.iplantcollaborative.org/browse/iplant/home/shared/onekp\_pilot]
- iPlant User Registration.** [http://user.iplantcollaborative.org]
- iPlant Learning Center.** [http://www.iplantcollaborative.org/learning-center/all-tutorials]
- iPlant Discovery Environment.** [http://de.iplantcollaborative.org]
- Using the iDrop Desktop.** [https://pods.iplantcollaborative.org/wiki/display/DS/Using+iDrop+Desktop]
- Matasci N, McKay SJ: **Phylogenetic analysis with the iPlant discovery environment.** *Curr Protoc Bioinformatics* 2013, **6**:Unit 6.13.
- Newick Trees Format.** [http://evolution.genetics.washington.edu/phyliip/newicktree.html]
- iRODS Data Management Software** [http://irods.org]
- Using iCommands (Unix).** [https://pods.iplantcollaborative.org/wiki/display/DS/Using+iCommands]
- Minie M, Chopra G, Sethi G, Horst J, White G, Roy A, Hatti K, Samudrala R: **CANDO and the infinite drug discovery frontier.** *Drug Discov Today* 2014, **19**:1353–1363.
- BioGRID Interactions.** [http://thebiogrid.org]
- Kyoto Encyclopedia of Genes and Genomes (KEGG).** [http://www.genome.jp/kegg]
- 1KP Protein-Protein Interactions Mapped to Metabolic Pathways.** [http://proinfo.org/1kp]
- 1KP Capstone Objective.** [https://pods.iplantcollaborative.org/wiki/display/iptol/OneKP+Capstone+Wiki]
- 1KP Papers in Progress.** [https://pods.iplantcollaborative.org/wiki/display/iptol/OneKP+companion+papers]
- Sayou C, Monniaux M, Nanao MH, Moyroud E, Brockington SF, Thévenon E, Chahtane H, Warthmann N, Melkonian M, Zhang Y, Wong GK, Weigel D, Parcy F, Dumas R: **A promiscuous intermediate underlies the evolution of LEAFY DNA binding specificity.** *Science* 2014, **343**:645–648.
- Klapoetke NC, Murata Y, Kim SS, Pulver SR, Birdsey-Benson A, Cho YK, Morimoto TK, Chuong AS, Carpenter EJ, Tian Z, Wang J, Xie Y, Yan Z, Zhang Y, Chow BY, Surek B, Melkonian M, Jayaraman V, Constantine-Paton M, Wong GK, Boyden ES: **Independent optical excitation of distinct neural populations.** *Nat Methods* 2014, **11**:338–346.

doi:10.1186/2047-217X-3-17

**Cite this article as:** Matasci et al.: Data access for the 1,000 Plants (1KP) project. *GigaScience* 2014 **3**:17.